



Railway Scheduling with Hierarchical Reinforcement Learning on a Graph

Thomas Dubach Florian Fuchs Francesco Corman

Institute for Transport Planning and Systems, ETH Zürich, Switzerland

1 Introduction

Railway scheduling is a complex task that involves solving recurring instances of constrained optimization problems. Traditionally, each instance is treated as an isolated problem, leading to inefficiencies in computational effort and scalability. This research explores reinforcement learning for optimizing railway timetables based on a graph representation of the scheduling problem.

The graph has nodes which represent operations which have a minimum running time δ , lower and upper bounds for the scheduled time t and a resource which is occupied by the train between $[t, t + \delta]$. The links represent dependencies where exactly one outgoing node must follow immediately after an operation finishes. Thus multiple outgoing links represent alternative routes for a train.

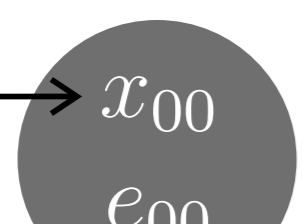
2 Graph embedding & reinforcement learning

We represent the graph state for the RL agent using a graph neural network. A feature vector x_i collects the information on each node i . Via message passing the features on each node are aggregated with the embeddings from its outgoing neighbors (i.e. messages are propagated bottom-up).

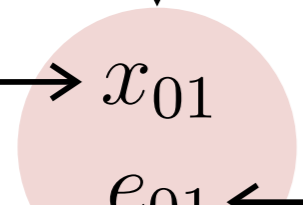
The RL agent uses scoring functions to decide which train to schedule next, how many of its nodes to schedule, and which route to take at junctions. These functions take the node embeddings as input and output probability distributions over the respective action spaces. By sampling from these distributions, the agent thus selects a subset of nodes to schedule. Standard RL algorithms such as Actor-Critic update both the scoring functions and the embeddings during training.

features train 0

$t = 1$
 $\delta \geq 2$
 r_1



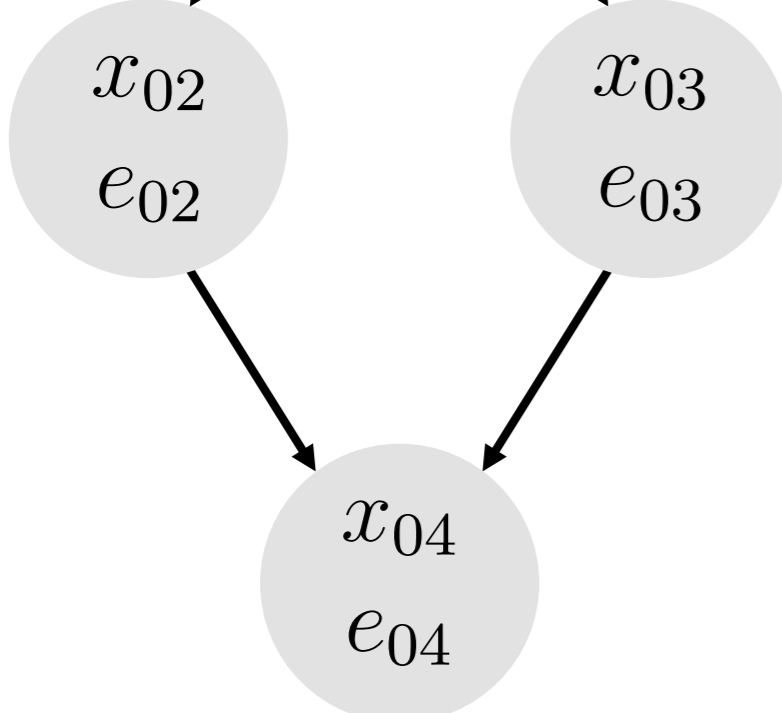
features
 $0 \geq t \geq 5$
 $\delta \geq 4$
 r_1



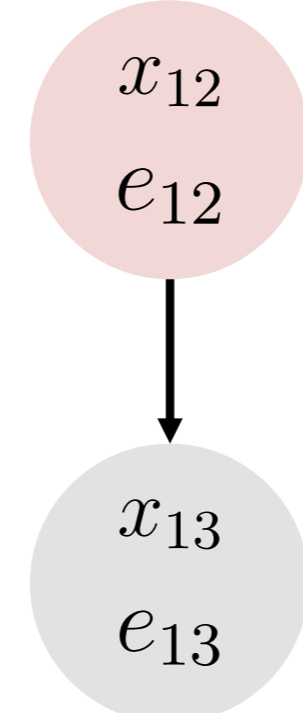
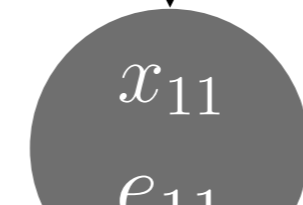
embedding

message passing

$$e_i = g \left(\sum_{j \in \mathcal{N}^{\text{out}}(i)} f(e_j) \right) + h(x_i)$$



train 1



● scheduled node

● next unscheduled node

● unscheduled node

t scheduled time

δ minimal running time

r_1 resource occupation

f, g, h neural networks

\mathcal{N}^{out} outgoing neighbors

s scoring neural networks

train scoring $s_{\text{train}}(e_{01}), s_{\text{train}}(e_{12}) \xrightarrow{\text{softmax}} p_0, p_1 \xrightarrow{\text{sample}} \text{schedule train 0}$

number scoring $s_{\text{num}}(e_{01}, 1), s_{\text{num}}(e_{01}, 2), \dots, s_{\text{num}}(e_{01}, l_{\text{max}}) \xrightarrow{\text{softmax}} p_1, p_2, \dots, p_{l_{\text{max}}} \xrightarrow{\text{sample}} \text{schedule 2 nodes}$

number scoring $s_{\text{route}}(e_{02}), s_{\text{route}}(e_{03}) \xrightarrow{\text{softmax}} p_2, p_3 \xrightarrow{\text{sample}} \text{schedule node 03}$

3 Reward based on feasibility

Given a set of scheduled nodes and implied precedences between trains, a timetable feasibility solver finds a feasible solution for the remaining unscheduled part. This feasible solution admits an objective value. As reward for the RL agent, we use the change of the objective value with respect to the solution found in the previous step.

If the feasibility solver does not find any solution, it returns a set of nodes causing the infeasibility. These can be traced back to decisions taken by the agent which are penalized.

4 Open questions

- How to represent precedences between trains in the graph?
- How to integrate the feasible solution in the environment state?

5 Selected references

- Mao, H., Schwarzkopf, M., Venkatakrisnan, S. B., Meng, Z., & Alizadeh, M. (2019). Learning scheduling algorithms for data processing clusters.
- Grinsztajn, N. (2023). Reinforcement learning for combinatorial optimization: Leveraging uncertainty, structure and priors (Université de Lille).